Atmospheric
Measurement
Techniques
Discussions

# Estimation of PM$_{2.5}$ Concentration in China Using Linear Hybrid Machine Learning Model

Zhihao Song[1], Bin Chen[1], Yue Huang[1], Li Dong [1], Tingting Yang [2]

[1]Atmospheric Science College of Lanzhou University, Lanzhou 730000, China

[2]Gansu Seed General Station, Lanzhou 730030, China

*Correspondence to*: Bin Chen (chenbin@lzu.edu.cn)

**Abstract.** The satellite remote-sensing aerosol optical depth (AOD) and meteorological elements were employed to invert PM$_{2.5}$ in order to control air pollution more effectively. This paper proposes a restricted gradient-descent linear hybrid machine learning model (RGD–LHMLM) by integrating a random forest (RF), a gradient boosting regression tree (GBRT), and a deep neural network (DNN) to estimate the concentration of PM$_{2.5}$ in China in 2019. The research data included Himawari-8 AOD with high spatiotemporal resolution, ERA-5 meteorological data, and geographic information. The results showed that, in the hybrid model developed by linear fitting, the DNN accounted for the largest proportion, whereas the weight coefficient was 0.62. The $R^2$ values of RF, GBRT, and DNN were reported 0.79, 0.81, and 0.8, respectively. Preferably, the generalization ability of the mixed model was better than that of each sub-model, and $R^2$ reached 0.84, whereas RMSE and MAE were reported 12.92 µg/m$^3$ and 8.01 µg/m$^3$, respectively. For the RGD-LHMLM, $R^2$ was above 0.7 in more than 70% of the sites, whereas RMSE and MAE were below 20 µg/m$^3$ and 15 µg/m$^3$, respectively, in more than 70% of the sites due to the correlation coefficient having seasonal difference between the meteorological factor and PM$_{2.5}$. Furthermore, the hybrid model performed best in winter (mean $R^2$ was 0.84) and worst in summer (mean $R^2$ was 0.71). The spatiotemporal distribution characteristics of PM$_{2.5}$ in China were then estimated and analyzed. According to the results, there was severe pollution in winter with an average concentration of PM$_{2.5}$ being reported 62.10 µg/m$^3$. However, there was slight pollution in summer with an average concentration of PM$_{2.5}$ being reported 47.39 µg/m$^3$. The findings also indicate that North China and East China are more polluted than other areas and that their average annual concentration of PM$_{2.5}$ was reported 82.68 µg/m$^3$. Moreover, there was relatively low pollution in Inner Mongolia, Qinghai, and Tibet, for their average PM$_{2.5}$ concentrations were reported below 40 µg/m$^3$.

1    **1 Background**

2       In recent years, pollutants have been discharged increasingly in China where air pollution is

3    becoming worse than ever before due to rapid urbanization and industrialization (Wang et al., 2019a).

4    The fine particulate matter (PM2.5) with a diameter below 2.5μm is the main component of air pollutants

5    having considerable impacts on human health, atmospheric visibility, and climate change (Gao et al.,

6    2015;Pan et al., 2018;Pun et al., 2017). The global concern about PM2.5 has increased significantly since

7    it was listed as a top carcinogen (Apte et al., 2015;Lim et al., 2020). Currently, ground monitoring is the

8    most efficient method of measuring $PM_{2.5}$ (Yang et al., 2018). However, monitoring stations are not

9    evenly distributed due to terrain and construction costs; therefore, it is difficult to obtain a wide range of

10   accurate $PM_{2.5}$ concentration data (Han et al., 2015). To solve the problem, the method of estimating

11   $PM_{2.5}$ with satellite remote-sensing was developed. Satellite remote-sensing is characterized by a wide

12   coverage and high resolution (Hoff and Christopher, 2009;Xu et al., 2021). There is also a high

13   correlation between AOD, obtained from satellite remote sensing inversion, and $PM_{2.5}$; therefore, AOD

14   is a very effective method of monitoring the spatiotemporal concentration characteristics of $PM_{2.5}$.

15      After Engel-Cox et al. (2004) proposed using satellite AOD to estimate $PM_{2.5}$ concentration, several

16   studies are reported in the literature to address this theory. Based on the regression model, Liu et al. (2005)

17   introduced AOD, boundary layer height, relative humidity, and geographical parameters as the main

18   controlling factors to estimate $PM_{2.5}$ in the eastern part of the United States, and the verification

19   coefficient $R^2$ obtained was 0.46. Tian and Chen (2010) used AOD, $PM_{2.5}$, and meteorological parameters

20   in Southern Ontario, Canada, to establish a semi-empirical model to predict $PM_{2.5}$ concentration per hour,

21   and the verification coefficient $R^2$ obtained in rural and urban areas was 0.7 and 0.64, respectively. Hu et

22   al. (2013) proposed a geography weighted regression model to estimate the surface $PM_{2.5}$ concentration

23   in southeastern America by combining AOD, meteorological parameters, and land use information. Their

24   model average $R^2$ was 0.6. Lee et al. (2012) believed that the satellite remote sensing AOD data would

25   be interfered by clouds and snow and ice, and the reliability of the data was questionable. They proposed

26   a mixed model based on AOD calibration to predict the ground $PM_{2.5}$ concentration in New England,

27   USA, and achieved good results ($R^2 = 0.83$). Combined with MODIS AOD and ground observation data,

28   Lv et al. (2017) estimated the daily surface $PM_{2.5}$ concentration in the Beijing-Tianjin-Hebei region and

29   improved the data resolution to 4 km. The data used in these early studies are AOD products obtained

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

1    from polar-orbit satellite sensors. The daily observation frequency is limited. Due to the influence of

2    cloud and ground reflection, the dynamic change information of $PM_{2.5}$ cannot be obtained. As a result,

3    geostationary satellite observations can be used to overcome the problem of low temporal resolution for

4    estimating surface $PM_{2.5}$ (Emili et al., 2010).

5        The Himawari-8 satellite commonly used in the Asia-Pacific region is a geostationary satellite

6    launched by the Japan Meteorological Agency in 2014. The observation frequency is 10 minutes, and the

7    observation results can characterize the aerosol and provide AOD data with a resolution of 5 km (Bessho

8    et al., 2016;Yumimoto et al., 2016). Due to its excellent performance, some scholars use Himawari-8

9    data to estimate ground $PM_{2.5}$.Wang et al. (2017) proposed an improved linear model, introduced AOD,

10    meteorological parameters, geographic information to estimate $PM_{2.5}$ in the Beijing-Tianjin-Hebei region,

11    and the verification coefficient R² was 0.86. Zhang et al. (2019b) used Himawari-8 hourly AOD product

12    to estimate ground $PM_{2.5}$ in China's four major urban agglomerations. The results showed significant

13    diurnal, seasonal, and spatial changes and improved the temporal resolution of estimating $PM_{2.5}$

14    concentration to the hourly level.

15        As research into ground-based $PM_{2.5}$ estimation deepens, traditional linear or nonlinear models

16    cannot meet the requirements of large-scale estimation and are gradually being replaced by machine

17    learning algorithms with strong nonlinear fitting ability. Liu et al. (2018) combined Kriging interpolation

18    and random forest algorithm to obtain the concentration of high-resolution ground $PM_{2.5}$ in the United

19    States. To demonstrate the accuracy and superiority of the proposed method, the results were compared

20    with the $PM_{2.5}$ concentration in ground measurement stations. Chen et al. (2019) stacked and predicted

21    $PM_{2.5}$ concentration based on a variety of machine learning algorithms, discussed the influence of

22    meteorological factors on $PM_{2.5}$ and achieved an $R^2 = 0.85$. Li et al. (2017a) established a GRNN model

23    for the whole of China to estimate $PM_{2.5}$ concentration, and the results demonstrated that the performance

24    of the deep learning model was better than that of the traditional linear model.

25        A large number of existing studies in the broader literature have examined the estimation of ground

26    $PM_{2.5}$ concentrations using satellite remote sensing AOD. However, the performance of $PM_{2.5}$ estimation

27    models established in the existing studies varies greatly and the performance of the models is not stable

28    in different seasons and regions. To overcome this limitation, in this paper, a linear hybrid machine

29    learning model (RGD-LHMLM) based on random forest (RF), gradient lifting regression tree (GBRT),

30    and deep neural network (DNN) is proposed to estimate ground $PM_{2.5}$ concentration. The model
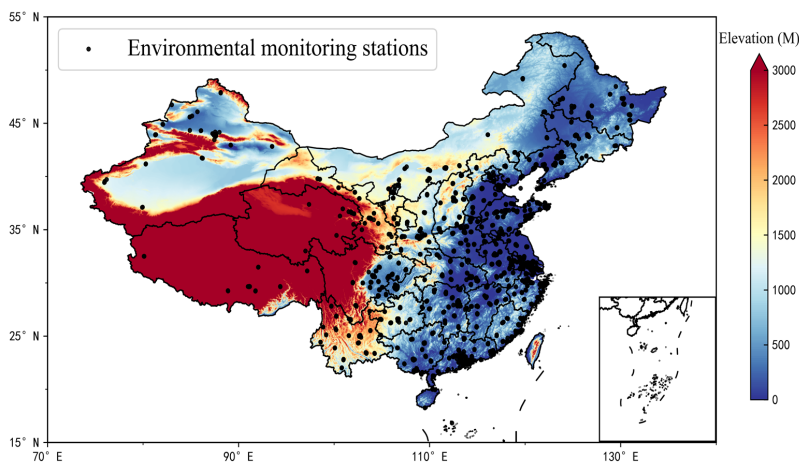
Atmospheric
Measurement
Techniques
Discussions

1    performance is evaluated from time and space to analyze its causes. Finally, spatiotemporal distribution

2    of PM$_{2.5}$ concentration in China in 2019 is obtained.

3    **2 Data**

4    **2.1 Ground PM$_{2.5}$ Monitoring Data**

5    PM$_{2.5}$ concentration data for 2019 used in this study are available from the China Environmental

6    Monitoring Center's Air Quality Real-Time Publication System. The system extracts hourly mean PM$_{2.5}$

7    data. By the end of 2019, China had 1641 monitoring stations built and in operation. Figure 1 shows the

8    spatial distribution of monitoring stations in China.



9
10    **Figure 1 Distribution diagram of Environmental monitoring stations in China (2019)**

11    **2.2 Satellite AOD Data**

12    The Himawari Imager (AHI) on the Himawari-8 satellite launched by the Japan Meteorological

13    Agency is a highly improved multi-wavelength imager. It adopts the whole disk observation method and

14    has 16 visible and infrared channels. It has the characteristics of fast imaging speed, flexible observation

15    area, and time. The Level-3-hour AOD product, released by the Japan Aerospace Space Agency (JAXA),

16    provides 500 nm AOD data with a spatial resolution of 5km during the day. In previous studies (Zang et

17    al., 2018), Himawari-8 AOD was compared with the AOD data of AERONET (Aerosol Robotic Network)

18    in China and achieved good performance. The AOD data used in this study is the Himawari-8 Level 3-

19    hour AOD data in 2019 obtained from the Himawari Monitor website of the Japan Meteorological

Atmospheric
Measurement
Techniques
Discussions
Open Access
EGU

1    Agency.

2    **2.3 Meteorological Data**

3    ERA-5 reanalysis data is an hourly collection of atmospheric and land-surface meteorological

4    elements since 1979 that the European Centre (ECMWF) has used its prediction model and data

5    assimilation system to "Reanalyse" archived observations. Data used in this paper include surface relative

6    humidity (RH, expressed as a percentage), air temperature at a height of 2 m (TM, expressed as K), Wind

7    speed (U10, V10, in m/s), surface pressure (SP, in Pa), boundary layer height (BLH, in m) and cumulative

8    precipitation (RAIN, in m) at 10 m above the ground. A series of studies has indicated that these

9    parameters can affect the concentration of $PM_{2.5}$ (Fang et al., 2016;Guo et al., 2017;Li et al., 2017b;Wang

10   et al., 2019b).

11   **2.4 Auxiliary Data**

12   The auxiliary data used in this study include high and low vegetation index (LH, LL),

13   ground elevation data (DEM), and population density data (PD). The high and low vegetation

14   index is derived from ERA5 reanalysis data, which respectively represent half of the total green

15   leaf area per unit level ground area of high and low vegetation type. The ground elevation data

16   are derived from SRTM-3 measurements jointly conducted by NASA and the Defense

17   Department's National Mapping Agency (NIMA), with a spatial resolution of 90 m. The

18   population data come from the 2015 United Nations Adjust Population Density data provided

19   by NASA's Center for Socio-Economic Data and Applications (SEDAC), which is based on

20   national censuses and adjusted for relative spatial distribution.

21   **3 Method**

22   **3.1 Random Forest**

23   Random Forest (RF) is built based on the combination of the Bagging algorithm and decision tree,

24   which is an extended variant of the parallel ensemble learning method (Stafoggia et al., 2019). To

25   construct a large number of decision trees, the random forest model takes multiple samples of the sample

26   data. In the decision tree, the nodes are divided into sub-nodes by using the randomly selected optimal

27   features until all the training samples of the node belong to the same class. Finally, all the decision trees

1    are merged to form the random forest. This method has proved to be effective in regression and

2    classification problems and is one of the most well-known Machine learning algorithms used in many

3    different fields (Yesilkanat, 2020).

4    **3.2 Gradient Boosted Regression Trees**

5        Different from the random forest, Gradient Boosting Regression Tree (GBRT) is based on Boosting

6    algorithm and decision tree. The basic principle of GBRT is to construct M different basic learners

7    through multiple iterations, and constantly add the weight of the learners with a small error probability,

8    to eventually generate a strong learner (Johnson et al., 2018). The core of this method is that after each

9    iteration, a learner will be built in the direction of residual reduction (gradient direction) to make the

10    residual decrease in the gradient direction (Schonlau, 2005). The basic learner of GBRT is the regression

11    tree in the decision tree. During the prediction, a predicted value is calculated according to the model

12    obtained. The minimum square root error is used to select the optimal feature to split the dataset, and the

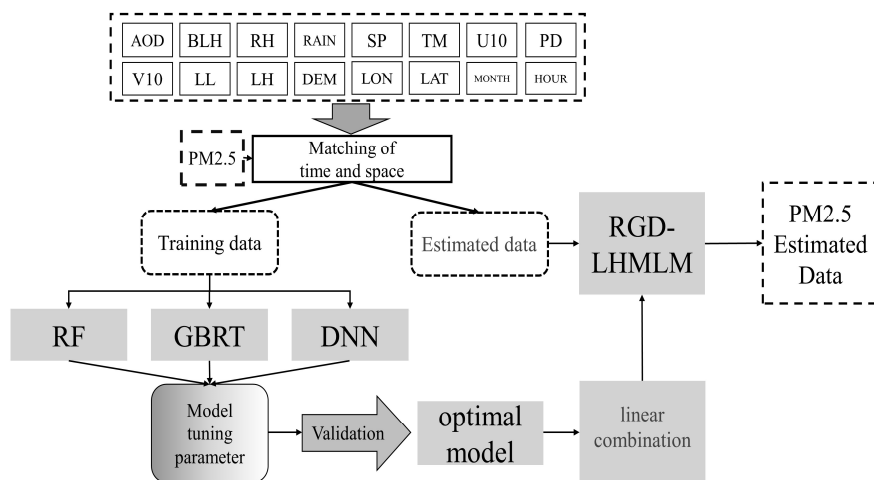13    average value of the child node is then taken as the predicted value.

14    **3.3 Deep Neural Networks**

15        Deep Neural Networks (DNN) is a supervised learning technique that uses a backpropagation

16    algorithm to minimize the loss function. It adjusts the parameters through an optimizer, and has high

17    computational power, making it ideal for solving classification and regression problems (Wang and Sun,

18    2019). The structure of DNN includes an input layer, an output layer, and several hidden layers. Each

19    layer takes the output of all nodes of the previous layer as the input, and this process requires activation

20    functions. Compared with other activation functions, the linear rectifying function (ReLU) has the

21    advantages of simple derivation, faster convergence, and higher efficiency. At the same time, among the

22    adaptive learning rate optimizers, the Adamx optimizer performs the best. It not only has the advantages

23    of Adam in determining the learning rate range and having stable parameters in each iteration but also

24    simplifies the method of defining the upper limit range of the learning rate and improves the iteration

25    efficiency (Diederik and Jimmy, 2015). Therefore, in this paper, we selected the Adamx optimizer and

26    ReLU activation function to train the DNN.

1 **3.4 Model Establishment and Verification**

2 After data processing, RF, GBRT, and DNN are used for modeling. To prevent model parameters

3 from being controlled by large or small range data and speed up the convergence rate of the model, the

4 data must be normalized before starting the training process. Finally, the three optimal sub-models are

5 linear combined to achieve the final mixed model. To verify the model performance, this paper uses the

6 "10-fold cross-validation" method (Adams et al., 2020). In this method, the data is split into 10 copies, 9

7 copies for training and 1 copy for verification; this process is repeated 10 times, and then the average of

8 the 10 predictions is computed as the final result. Finally, the predicted value and the measured value are

9 fitted linearly. At the same time, several indicators are used to evaluate the model, including the mean

10 absolute error (MAE, when the predicted value and the true value are exactly equal to 0, that is, perfect

11 model; The larger the error, the greater the value), the root mean square error (RMSE, when the predicted

12 value and the real value are completely consistent is equal to 0, that is, the perfect model; The larger the

13 error, the greater the value), the slope of the fitting equation and the determination coefficient $R^2$ (the

14 greater the value, the better the model fitting effect).



15
16 **Figure 2 Schematic diagram of model**

17 **4 Results and Discussion**

18 **4.1 Modeling Results**

19 According to the above steps, the mixed model RGD-LHMLM is obtained through modeling

Atmospheric
Measurement
Techniques
Discussions

1    verification, and is compared with RF, GBRT, and DNN. The fitting and verification accuracy results of
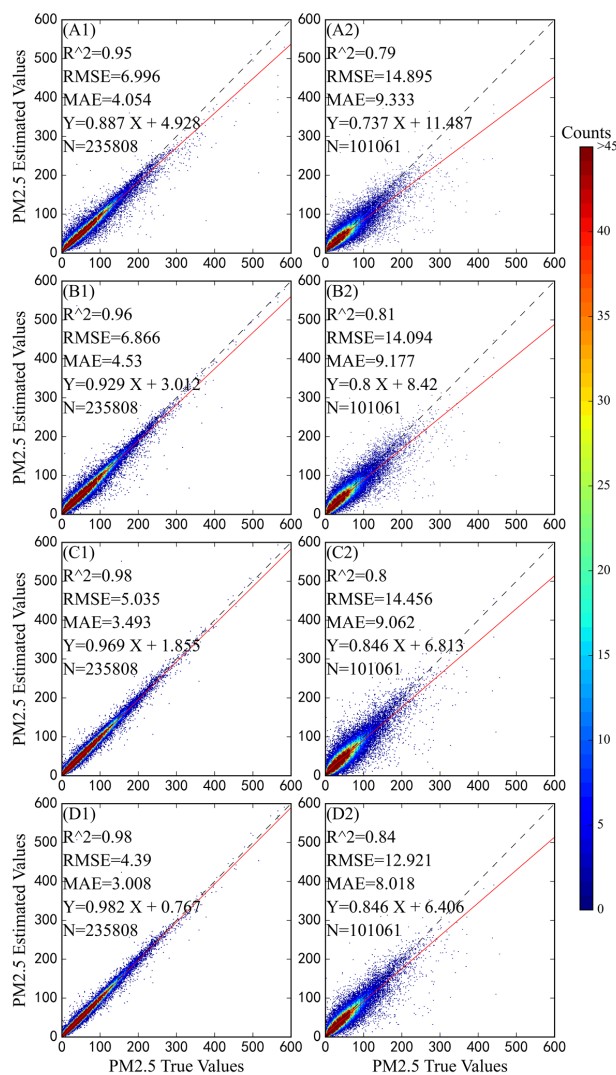
2    each model are shown in Table 1.

3

**Table 1 Comparison of model accuracy**

| Model | Fitting | | | Validation | | |
|---|---|---|---|---|---|---|
| | $R^2$ | RMSE | MAE | $R^2$ | RMSE | MAE |
| RF | 0.95 | 6.99 | 4.05 | 0.79 | 14.89 | 9.33 |
| GBRT | 0.96 | 6.87 | 4.52 | 0.81 | 14.09 | 9.18 |
| DNN | 0.97 | 5.03 | 3.49 | 0.80 | 14.45 | 9.06 |
| RGD-LHMLM | 0.98 | 4.39 | 3.00 | 0.84 | 12.92 | 8.01 |

4

5    The PM$_{2.5}$ inversion results of a single machine learning model show that DNN has the best

6    inversion performance, followed by GBRT, and RF has the worst performance. The expression of the

7    mixing model obtained after linear mixing is as follows:

8    $$PM_{2.5RGD-LHMLM} = 0.25PM_{2.5RF} + 0.17PM_{2.5GBRT} + 0.62PM_{2.5DNN} - 2.13 \qquad (1)$$

9    The weight coefficient of DNN in the mixed model was the largest (0.62). The $R^2$ of RGD-LHMLM in

10   the training set was 0.98, and the RMSE was only 4.39 μg/m$^3$, indicating that the model had an excellent

11   data fitting effect. Meanwhile, the generalization ability of the mixed model is also good, with $R^2$ of 0.84

12   and RMSE of 12.92 μg/m$^3$ on the validation data set. Compared with RF, GBRT, and DNN, the inversion

13   performance of RGD-LHMLM is significantly improved. In other words, the combination of multiple

14   models can improve the robustness and generalization ability of the model (Wolpert, 1992). The linear

15   fitting equation coefficients between the predicted and measured values in the training set and the

16   verification set were 0.98 and 0.84, respectively, indicating that the prediction accuracy of the model

17   reached a high level. The fitting curve between the model predicted value and the real value is shown in

18   Figure 3. The RGD-LHMLM model has the smallest degree of data dispersion, and the slope of the fitting

19   line reaches 0.84, indicating that 84% of the prediction results are accurate, higher than the three sub-

20   models.

1

2 **Figure 3 Accuracy of model Fitting and Validation (A: RF, B: GBRT, C: DNN, D: RGD-LHMLM)**

3 **4.2 Model Performance Analysis**

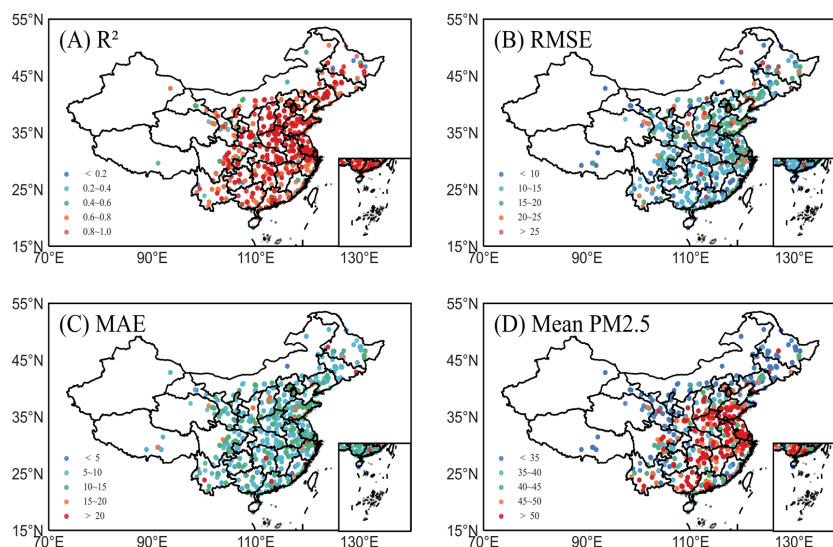4 **4.2.1 Performance Analysis of Monitoring Station Model**

5     The spatial performance of the model was analyzed by measuring $R^2$, RMSE, and MAE at the

6 monitoring stations. According to Figure 4, there are regional differences in the inversion performance

7 of RGD-LHMLM. At all monitoring stations, the average $R^2$ was reported 0.74, and $R^2$ was above 0.7 at

8 more than 70% of the stations, especially in the densely populated and industrially developed areas. The

1   model prediction accuracy was reported low ($R^2<0.6$) in Xinjiang, Tibet, Qinghai, Western Sichuan, and

2   a few other areas of Northeast China. The mean values of RMSE and MAE were reported 11.4 μg/m$^3$

3   and 8.01 μg/m$^3$, respectively. In fact, the mean values of RMSE and MAE were below 20 μg/m$^3$ and 15

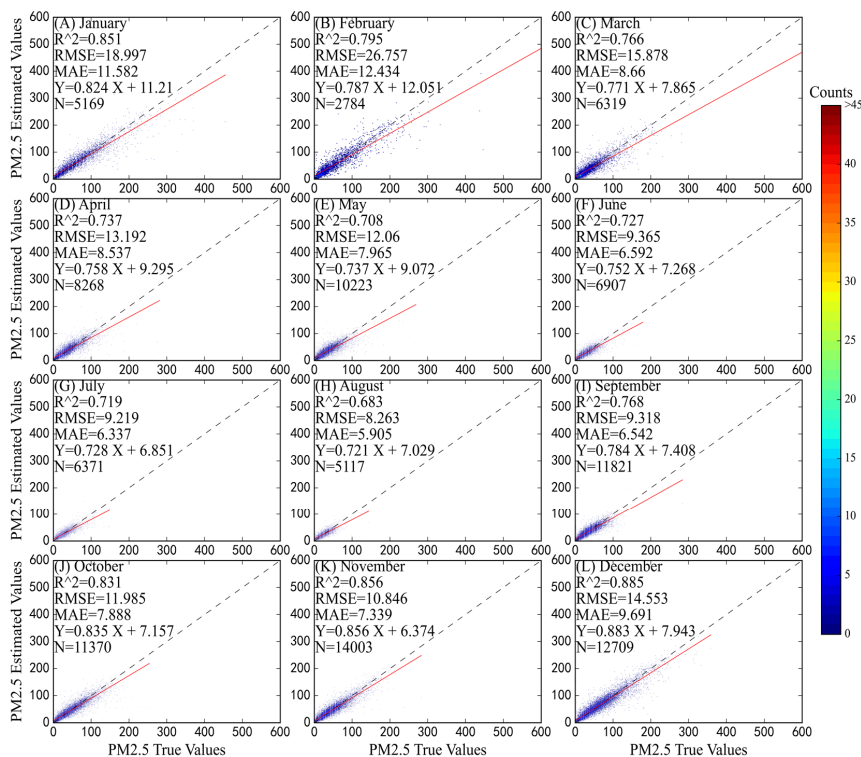4   μg/m$^3$ in more than 95% of stations, something showed a low estimation error.



5
6   **Figure 4 Model precision parameters (A)$R^2$, (B)RMSE, (C)MAE and (D)Mean PM$_{2.5}$ concentration site**

7   **distribution**

8   Based on the analysis of spatial differences in the RGD-LHMLM inversion performance, the

9   following deductions can be made. First, the environmental monitoring stations in the central and eastern

10  regions with better inversion performance were distributed densely, and there are large data available;

11  therefore, the model had a satisfactory training effect. Moreover, data matching was lower in the western

12  region than in other regions, something which resulted in model over-fitting and reduced accuracy

13  (Zhang et al., 2018). Second, some areas of western and northeastern China are covered by snow and the

14  Gobi Desert with high surface albedo. This reduces the accuracy of AOD obtained by satellite

15  observation and brings errors to model training. Finally, the Himawari-8 scanning range is limited, and

16  the satellite observation data obtained in Western China are limited in terms of quantity and accuracy. In

17  general, the RGD-LHMLM has a satisfactory spatial performance, especially in areas with high annual

18  average concentration of PM$_{2.5}$; therefore, it can leave a good inversion effect.

1    **4.2.2 Time-Scale Model Performance Analysis**

2        Figure 5 shows the inversion performance results of the hybrid model collected from January to

3    December 2019. The model performed the worst in summer months because $R^2$ was reported 0.73, 0.72,

4    and 0.68, respectively; however, RMSE and MAE were only 9.37, 9.22, 8.26 μg/m³ and 6.59, 6.34, and

5    5.91 μg/m³, respectively, due to the lower average concentration of $PM_{2.5}$ in summer. Winter and autumn

6    models gained better performance results with an average $R^2$ over 0.8. However, in contrast to summer,

7    the estimation errors of these two seasons were relatively large, with average RMSE of 20.10 μg/m³ and

8    10.72 μg/m³ and average MAE of 11.20 μg/m³ and 7.25 μg/m³, respectively. The mean $R^2$ was 0.74,

9    whereas the mean RMSE and MAE were 13.71 μg/m³ and 8.39 μg/m³, respectively.
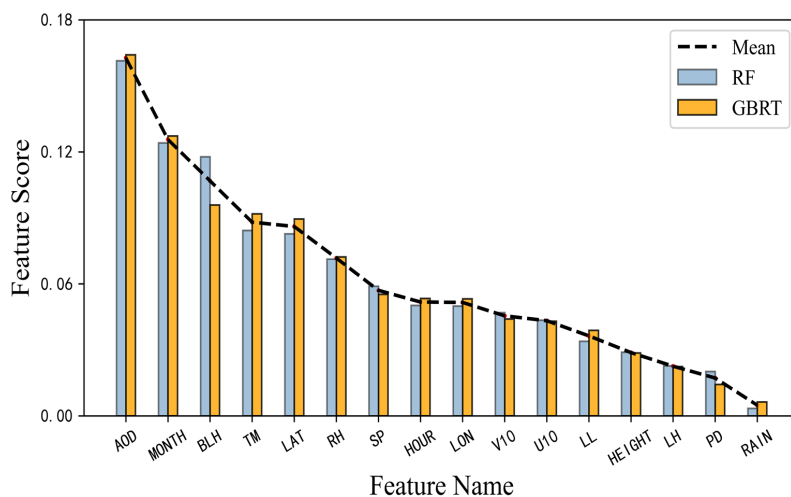


10

11        **Figure 5 Monthly model performance fitting scatter diagram in 2019**

12        The model performance differences were also analyzed to extract and rank the model features of

13    RF and GBRT based on the feature importance. The higher the feature importance, the greater the

14    contribution of factors to the model. Figure 6 shows that AOD, boundary layer height, 2 m surface

15    temperature, and relative humidity had the greatest effect on the mixed model performance out of all

1    variable characteristic parameters. Accordingly, AOD is greatly affected by the fine particulate matter

2    and is the main factor in the inversion of $PM_{2.5}$. Changes of the boundary layer height can affect the

3    diffusion ability of the atmosphere. If the boundary layer height is low, the accumulation of pollutants

4    will be caused. At the same time, the 2 m surface temperature has a great impact on the boundary layer

5    height (Miao et al., 2018). Finally, higher rates of atmospheric humidity can improve the fine particulate

6    matter accumulation.



7

8    **Figure 6 Importance of model features (represent the contribution of feature factors to the model)**

9    The correlation coefficients between the monthly mean values of important meteorological

10   parameters (AOD, BLH, TM and RH) and $R^2$ were also analyzed. According to the results, the correlation

11   coefficients between the meteorological parameters and $PM_{2.5}$ were lower in summer. Furthermore, there

12   are many rainy days and large cloud coverage, which is not conducive to satellite observation and

13   decreases the accuracy of AOD data in summer. Therefore, the summer model performance is poor. There

14   was a strong correlation between meteorological parameters and $PM_{2.5}$ in autumn. There were also

15   similar correlations between spring and winter; however, the winter model performed was better. The

16   reasons can be interpreted as below. The winter temperature and boundary layer height are low, whereas

17   the atmosphere is stable but not conducive to the diffusion of pollutants. Moreover, during the heating

18   period in winter, pollutant emissions soar greatly and result in a sharp rise in the concentration of $PM_{2.5}$.

19   The increased pollution in winter ensures the quality and quantity of data, thereby improving the model
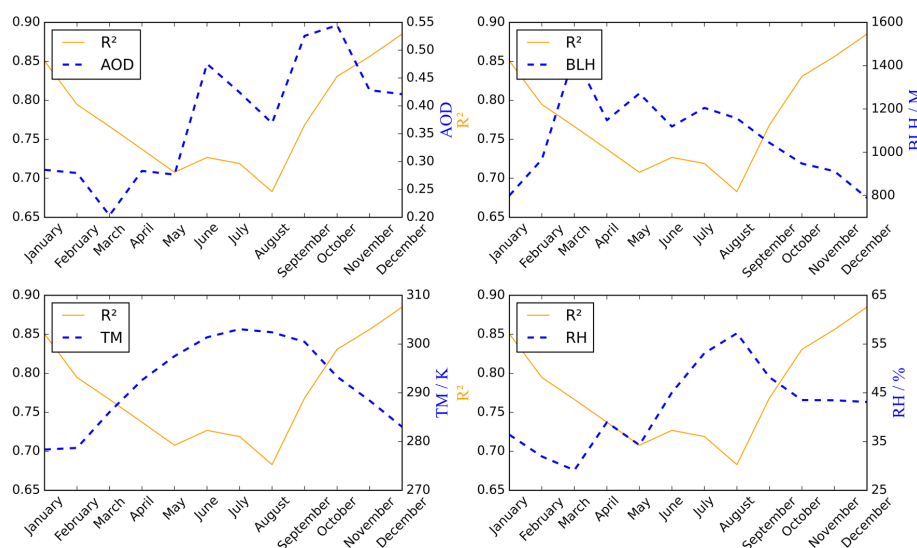
20   performance effectively.

1

2

**Table 2 Correlation coefficient between meteorological parameters with PM$_{2.5}$**

| Season | AOD | BLH | TM | RH |
|--------|-----|-----|----|----|
| Spring | 0.47 | -0.33 | 0.12 | 0.36 |
| Summer | 0.42 | -0.21 | 0.06 | 0.19 |
| Autumn | 0.38 | -0.29 | 0.24 | 0.41 |
| Winter | 0.44 | -0.33 | 0.12 | 0.35 |



3

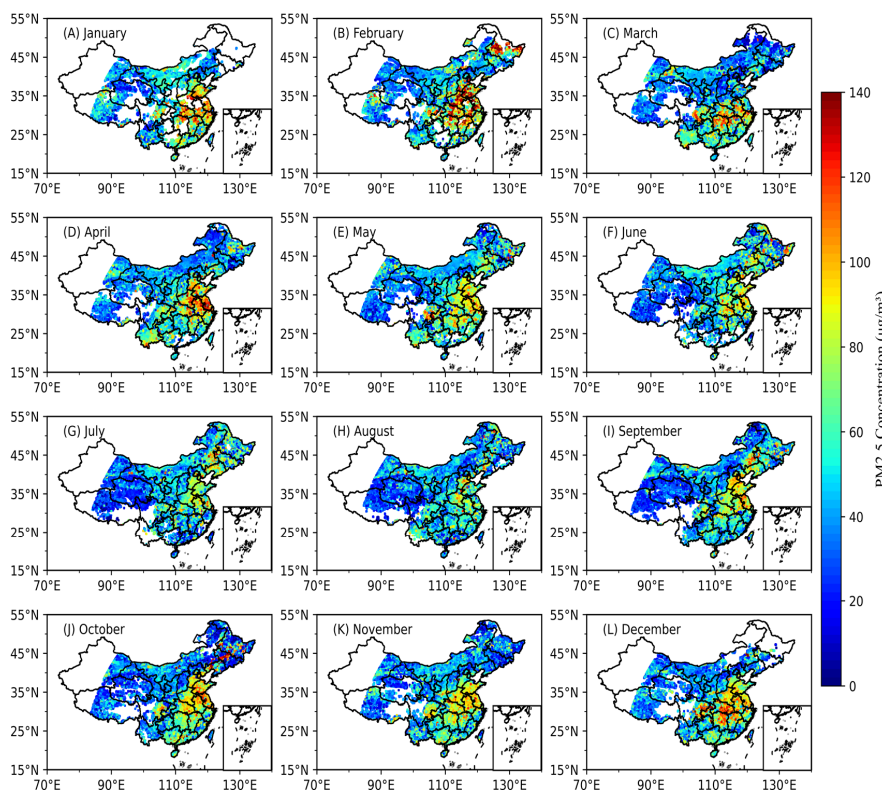4      **Figure 7 Variation trend of monthly average of meteorological parameters (AOD, BLH, TM, RH) and R$^2$**

5      **4.3 Temporal and Spatial Distribution Characteristics of PM$_{2.5}$ Concentration in China**

6      In terms of spatial distribution, Shandong, Henan, Jiangsu, Anhui, as well as parts of Hubei and

7      Hebei were the most polluted areas in China in 2019, with an annual average PM$_{2.5}$ concentration of

8      82.86 μg/m$^3$. On the one hand, these areas are economically developed and densely populated, resulting

9      in a large amount of pollutant emissions. On the other hand, the barrier of the peripheral mountains

10     (Taihang Mountains, Qinling Mountains and the Southern Hills) leads to the accumulation of pollutants

11     that are difficult to diffuse. Sichuan Basin is a rare area with a high PM$_{2.5}$ value due to its unique

12     topography (Zhang et al., 2019a), with an annual average PM$_{2.5}$ concentration of 64.69 μg/m$^3$. In addition,

13     Inner Mongolia, Qinghai, Tibet and other places, the pollution level is low, the average annual PM2.5

14     concentration is less than 40 μg/m$^3$.

15     PM$_{2.5}$ concentration in China varies significantly with the seasons. As shown in Figure 8, PM$_{2.5}$

16     concentration in winter is the highest, with an average value of 62.10μg/m$^3$. January 2019 was the most

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

1    polluted month in China, with the average $PM_{2.5}$ concentration reaching $63.58\mu g/m^3$. The average $PM_{2.5}$

2    concentration was 47.39 $\mu g/m^3$ in summer. The average concentration of $PM_{2.5}$ in spring and autumn was

3    $54.21\mu g/m^3$ and 52.26 $\mu g/m^3$, respectively, indicating similar levels of pollution.



4
5          **Figure 8 Monthly distribution of $PM_{2.5}$ concentration in China in 2019**

6    **5 Conclusion**

7        It is essential to collect the spatiotemporal evolution characteristics regarding the concentration of

8    $PM_{2.5}$ for air pollution prevention and containment. Based on the linear hybrid machine learning model,

9    this paper used the AOD data of Himawari-8 to invert the concentration of $PM_{2.5}$ in China and obtain its

10    distribution characteristics. The model performance and inversion results are analyzed and summarized

11    below:

12        (1) In the RGD-LHMLM obtained from linear fitting, the DNN accounted for the largest proportion

13    with a weight coefficient of 0.62. The $R^2$ of RGD-LHMLM was 0.84, whereas its generalization ability

14    was significantly better than that of a single model (DNN: 0.80; GBRT: 0.81; RF: 0.79). Moreover,

1    RMSE and MAE were 12.92 μg/m$^3$ and 8.01 μg/m$^3$, respectively.

2    (2) The RGD-LHMLM was spatially stable, with R$^2$>0.7 in more than 70% of sites as well as

3    RMSE<20 μg/m$^3$ and MAE<15μg/m$^3$ in more than 95% of sites. These sites are mainly located in densely

4    populated and industrially developed areas. The correlation difference between the inversion factor and

5    PM$_{2.5}$ in various seasons would lead to seasonal variations in the model performance. In addition, the

6    performance was the worst in summer with an average R$^2$ of 0.71; however, winter showed the best

7    performance with an average R$^2$ of 0.84.

8    (3) Changes in the spatiotemporal characteristics were obvious in the concentration of PM$_{2.5}$ in

9    China. In other words, North China and East China had the highest concentration of PM$_{2.5}$ with an

10   average annual concentration of 82.86 μg/m$^3$, whereas Inner Mongolia, Qinghai, Tibet, and other regions

11   had low pollution levels with an average annual concentration of PM$_{2.5}$ below 40 μg/m$^3$. In winter, the

12   concentration of PM$_{2.5}$ was higher with an average of 62.10 μg/m$^3$, whereas the pollution was lighter in

13   summer with an average concentration of PM$_{2.5}$ being reported 47.39 μg/m$^3$.

14   In conclusion, the RGD-LHMLM can accurately measure the concentration of PM$_{2.5}$ and perform

15   the seasonal evolution of pollutants. These results can help control the local pollution. This study also

16   indicated that integrating multiple Machine learning models improved the accuracy of fitting results

17   effectively. For more accurate pollutant data, such models can be employed to fit the PM$_{2.5}$ in the future

18   with more parameters closely related to PM$_{2.5}$. However, there are some vacant values in the results of

19   this study. There are also no data for some areas. Thus, other satellite data can be used in future studies

20   to solve this problem.

21   **Data availability**

22   Datasets related to this paper can be requested from the corresponding author (chenbin@lzu.edu.cn).

23   **Author contributions**

24   Chen proposed the content of the study. Song performed data processing, model building, result analysis,

25   and article writing. Huang, Dong and Yang checked the content of the article.

**Competing interests**

The authors declare that they have no conflict of interest.

**Acknowledgments**

**Financial support**

**References**

Adams, M. D., Massey, F., Chastko, K., and Cupini, C.: Spatial modelling of particulate matter air pollution sensor measurements collected by community scientists while cycling, land use regression with spatial cross-validation, and applications of machine learning for data correction,Atmos Environ, 230,https://doi.org/10.1016/j.atmosenv.2020.117479, 2020.

Apte, J. S., Marshall, J. D., Cohen, A. J., and Brauer, M.: Addressing Global Mortality from Ambient PM2.5,Environ Sci Technol, 49, 8057-8066,https://doi.org/10.1021/acs.est.5b01236, 2015.

Bessho, K., Date, K., Hayashi, M., Ikeda, A., Imai, T., Inoue, H., Kumagai, Y., Miyakawa, T., Murata, H., Ohno, T., Okuyama, A., Oyama, R., Sasaki, Y., Shimazu, Y., Shimoji, K., Sumida, Y., Suzuki, M., Taniguchi, H., Tsuchiyama, H., Uesawa, D., Yokota, H., and Yoshida, R.: An Introduction to Himawari-8/9-Japan's New-Generation Geostationary Meteorological Satellites,J Meteorol Soc Jpn, 94, 151-183,https://doi.org/10.2151/jmsj.2016-009, 2016.

Chen, J. P., Yin, J. H., Zang, L., Zhang, T. X., and Zhao, M. D.: Stacking machine learning model for estimating hourly PM2.5 in China based on Himawari 8 aerosol optical depth data,Sci Total Environ, 697,https://doi.org/10.1016/j.scitotenv.2019.134021, 2019.

Diederik, P. K., and Jimmy, B.: Adam: A Method for Stochastic Optimization, ICLR, 2015.

1  Emili, E., Popp, C., Petitta, M., Riffler, M., Wunderle, S., and Zebisch, M.: PM10 remote sensing from

2  geostationary SEVIRI and polar-orbiting MODIS sensors over the complex terrain of the European

3  Alpine region,Remote Sens Environ, 114, 2485-2499,https://doi.org/10.1016/j.rse.2010.05.024, 2010.

4  Engel-Cox, J. A., Holloman, C. H., Coutant, B. W., and Hoff, R. M.: Qualitative and quantitative

5  evaluation of MODIS satellite sensor data for regional and urban scale air quality,Atmos Environ, 38,

6  2495-2509,https://doi.org/10.1016/j.atmosenv.2004.01.039, 2004.

7  Fang, X., Zou, B., Liu, X. P., Sternberg, T., and Zhai, L.: Satellite-based ground PM2.5 estimation using

8  timely       structure     adaptive     modeling,Remote     Sens     Environ,     186,     152-

9  163,https://doi.org/10.1016/j.rse.2016.08.027, 2016.

10 Gao, M., Guttikunda, S. K., Carmichael, G. R., Wang, Y. S., Liu, Z. R., Stanier, C. O., Saide, P. E., and

11 Yu, M.: Health impacts and economic losses assessment of the 2013 severe haze event in Beijing area,Sci

12 Total Environ, 511, 553-561,https://doi.org/10.1016/j.scitotenv.2015.01.005, 2015.

13 Guo, J. P., Xia, F., Zhang, Y., Liu, H., Li, J., Lou, M. Y., He, J., Yan, Y., Wang, F., Min, M., and Zhai, P.

14 M.: Impact of diurnal variability and meteorological factors on the PM2.5 - AOD relationship:

15 Implications      for      PM2.5      remote      sensing,Environ      Pollut,      221,      94-

16 104,https://doi.org/10.1016/j.envpol.2016.11.043, 2017.

17 Han, Y., Wu, Y. H., Wang, T. J., Zhuang, B. L., Li, S., and Zhao, K.: Impacts of elevated-aerosol-layer

18 and aerosol type on the correlation of AOD and particulate matter with ground-based and satellite

19 measurements     in     Nanjing,     southeast     China,Sci     Total     Environ,     532,     195-

20 207,https://doi.org/10.1016/j.scitotenv.2015.05.136, 2015.

21 Hoff, R. M., and Christopher, S. A.: Remote Sensing of Particulate Pollution from Space: Have We

22 Reached the Promised Land?,J Air Waste Manage, 59, 645-675,https://doi.org/10.3155/1047-

23 3289.59.6.645, 2009.

24 Hu, X. F., Waller, L. A., Al-Hamdan, M. Z., Crosson, W. L., Estes, M. G., Estes, S. M., Quattrochi, D.

25 A., Sarnat, J. A., and Liu, Y.: Estimating ground-level PM2.5 concentrations in the southeastern US using

26 geographically weighted regression,Environ Res, 121, 1-10,https://doi.org/10.1016/j.envres.2012.11.003,

27 2013.

28 Johnson, N. E., Bonczak, B., and Kontokosta, C. E.: Using a gradient boosting model to improve the

29 performance of low-cost aerosol monitors in a dense, heterogeneous urban environment,Atmos Environ,

30 184, 9-16,https://doi.org/10.1016/j.atmosenv.2018.04.019, 2018.

1  Lee, H. J., Coull, B. A., Bell, M. L., and Koutrakis, P.: Use of satellite-based aerosol optical depth and

2  spatial clustering to predict ambient PM2.5 concentrations,Environ Res, 118, 8-

3  15,https://doi.org/10.1016/j.envres.2012.06.011, 2012.

4  Li, T. W., Shen, H. F., Zeng, C., Yuan, Q. Q., and Zhang, L. P.: Point-surface fusion of station

5  measurements and satellite observations for mapping PM2.5 distribution in China: Methods and

6  assessment,Atmos Environ, 152, 477-489,https://doi.org/10.1016/j.atmosenv.2017.01.004, 2017a.

7  Li, T. W., Shen, H. F., Yuan, Q. Q., Zhang, X. C., and Zhang, L. P.: Estimating Ground-Level PM2.5 by

8  Fusing Satellite and Station Observations: A Geo-Intelligent Deep Learning Approach,Geophys Res Lett,

9  44, 11985-11993,https://doi.org/10.1002/2017gl075710, 2017b.

10  Lim, C. H., Ryu, J., Choi, Y., Jeon, S. W., and Lee, W. K.: Understanding global PM2.5 concentrations

11  and their drivers in recent decades (1998-2016),Environ Int,

12  144,https://doi.org/10.1016/j.envint.2020.106011, 2020.

13  Liu, Y., Sarnat, J. A., Kilaru, A., Jacob, D. J., and Koutrakis, P.: Estimating ground-level PM2.5 in the

14  eastern united states using satellite remote sensing,Environ Sci Technol, 39, 3269-

15  3278,https://doi.org/10.1021/es049352m, 2005.

16  Liu, Y., Cao, G. F., Zhao, N. Z., Mulligan, K., and Ye, X. Y.: Improve ground-level PM2.5 concentration

17  mapping using a random forests-based geostatistical approach,Environ Pollut, 235, 272-

18  282,https://doi.org/10.1016/j.envpol.2017.12.070, 2018.

19  Lv, B. L., Hu, Y. T., Chang, H. H., Russell, A. G., Cai, J., Xu, B., and Bai, Y. Q.: Daily estimation of

20  ground-level PM2.5 concentrations at 4 km resolution over Beijing-Tianjin-Hebei by fusing MODIS

21  AOD and ground observations,Sci Total Environ, 580, 235-

22  244,https://doi.org/10.1016/j.scitotenv.2016.12.049, 2017.

23  Miao, Y. C., Liu, S. H., Guo, J. P., Huang, S. X., Yan, Y., and Lou, M. Y.: Unraveling the relationships

24  between boundary layer height and PM2.5 pollution in China based on four-year radiosonde

25  measurements,Environ Pollut, 243, 1186-1195,https://doi.org/10.1016/j.envpol.2018.09.070, 2018.

26  Pan, Z. X., Mao, F. Y., Wang, W., Zhu, B., Lu, X., and Gong, W.: Impacts of 3D Aerosol, Cloud, and

27  Water Vapor Variations on the Recent Brightening during the South Asian Monsoon Season,Remote

28  Sens-Basel, 10,https://doi.org/10.3390/rs10040651, 2018.

29  Pun, V. C., Kazemiparkouhi, F., Manjourides, J., and Suh, H. H.: Long-Term PM2.5 Exposure and

30  Respiratory, Cancer, and Cardiovascular Mortality in Older US Adults,Am J Epidemiol, 186, 961-

969,https://doi.org/10.1093/aje/kwx166, 2017.

Schonlau, M.: Boosted regression (boosting): An introductory tutorial and a Stata plugin,Stata J, 5, 330-354,https://doi.org/10.1177/1536867x0500500304, 2005.

Stafoggia, M., Bellander, T., Bucci, S., Davoli, M., de Hoogh, K., de'Donato, F., Gariazzo, C., Lyapustin, A., Michelozzi, P., Renzi, M., Scortichini, M., Shtein, A., Viegi, G., Kloog, I., and Schwartz, J.: Estimation of daily PM10 and PM2.5 concentrations in Italy, 2013-2015, using a spatiotemporal land-use random-forest model,Environ Int, 124, 170-179,https://doi.org/10.1016/j.envint.2019.01.016, 2019.

Tian, J., and Chen, D. M.: A semi-empirical model for predicting hourly ground-level fine particulate matter (PM2.5) concentration in southern Ontario from satellite remote sensing and ground-based meteorological measurements,Remote Sens Environ, 114, 221-229,https://doi.org/10.1016/j.rse.2009.09.011, 2010.

Wang, W., Mao, F. Y., Du, L., Pan, Z. X., Gong, W., and Fang, S. H.: Deriving Hourly PM2.5 Concentrations from Himawari-8 AODs over Beijing-Tianjin-Hebei in China,Remote Sens-Basel, 9,https://doi.org/10.3390/rs9080858, 2017.

Wang, X. H., Zhong, S. Y., Bian, X. D., and Yu, L. J.: Impact of 2015-2016 El Nino and 2017-2018 La Nina on PM2.5 concentrations across China,Atmos Environ, 208, 61-73,https://doi.org/10.1016/j.atmosenv.2019.03.035, 2019a.

Wang, X. P., and Sun, W. B.: Meteorological parameters and gaseous pollutant concentrations as predictors of daily continuous PM2.5 concentrations using deep neural network in Beijing-Tianjin-Hebei, China,Atmos Environ, 211, 128-137,https://doi.org/10.1016/j.atmosenv.2019.05.004, 2019.

Wang, X. Q., Wei, W., Cheng, S. Y., Yao, S., Zhang, H. Y., and Zhang, C.: Characteristics of PM2.5 and SNA components and meteorological factors impact on air pollution through 2013-2017 in Beijing, China,Atmospheric Pollution Research, 10, 1976-1984,https://doi.org/10.1016/j.apr.2019.09.004, 2019b.

Wolpert, D. H.: Stacked Generalization,Neural Networks, 5, 241-259,https://doi.org/10.1016/S0893-6080(05)80023-1, 1992.

Xu, J. H., Lindqvist, H., Liu , Q. F., Wang, K., and Wang, L.: Estimating the spatial and temporal variability of the ground-level NO2 concentration in China during 2005–2019 based on satellite remote sensing,,Atmospheric Pollution Research, 12, 57-67,https://doi.org/https://doi.org/10.1016/j.apr.2020.10.008, 2021.

Yang, X. C., Jiang, L., Zhao, W. J., Xiong, Q. L., Zhao, W. H., and Yan, X.: Comparison of Ground-

Atmospheric
Measurement
Techniques
Discussions

Open Access

Based PM2.5 and PM10 Concentrations in China, India, and the US,Int J Env Res Pub He, 15,https://doi.org/10.3390/ijerph15071382, 2018.

Yesilkanat, C. M.: Spatio-temporal estimation of the daily cases of COVID-19 in worldwide using random forest machine learning algorithm,Chaos Soliton Fract, 140,https://doi.org/10.1016/j.chaos.2020.110210, 2020.

Yumimoto, K., Nagao, T. M., Kikuchi, M., Sekiyama, T. T., Murakami, H., Tanaka, T. Y., Ogi, A., Irie, H., Khatri, P., Okumura, H., Arai, K., Morino, I., Uchino, O., and Maki, T.: Aerosol data assimilation using data from Himawari-8, a next-generation geostationary meteorological satellite,Geophys Res Lett, 43, 5886-5894,https://doi.org/10.1002/2016gl069298, 2016.

Zang, L., Mao, F. Y., Guo, J. P., Gong, W., Wang, W., and Pan, Z. X.: Estimating hourly PM1 concentrations from Himawari-8 aerosol optical depth in China,Environ Pollut, 241, 654-663,https://doi.org/10.1016/j.envpol.2018.05.100, 2018.

Zhang, L., Guo, X. M., Zhao, T. L., Gong, S. L., Xu, X. D., Li, Y. Q., Luo, L., Gui, K., Wang, H. L., Zheng, Y., and Yin, X. F.: A modelling study of the terrain effects on haze pollution in the Sichuan Basin,Atmos Environ, 196, 77-85,https://doi.org/10.1016/j.atmosenv.2018.10.007, 2019a.

Zhang, T. H., Zhu, Z. M., Gong, W., Zhu, Z. R., Sun, K., Wang, L. C., Huang, Y. S., Mao, F. Y., Shen, H. F., Li, Z. W., and Xu, K.: Estimation of ultrahigh resolution PM2.5 concentrations in urban areas using 160 m Gaofen-1 AOD retrievals,Remote Sens Environ, 216, 91-104,https://doi.org/10.1016/j.rse.2018.06.030, 2018.

Zhang, T. X., Zang, L., Wan, Y. C., Wang, W., and Zhang, Y.: Ground-level PM2.5 estimation over urban agglomerations in China with high spatiotemporal resolution based on Himawari-8,Sci Total Environ, 676, 535-544,https://doi.org/10.1016/j.scitotenv.2019.04.299, 2019b.